

Clustering with Spectral Norm and the k -means Algorithm

Amit Kumar *

Dept. of Computer Science and Engg.

IIT Delhi, New Delhi

email : amitk@cse.iitd.ac.in

Ravindran Kannan

Microsoft Research India Lab.

Bangalore

email : kannan@microsoft.com

April 13, 2010

Abstract

There has been much progress on efficient algorithms for clustering data points generated by a mixture of k probability distributions under the assumption that the means of the distributions are well-separated, i.e., the distance between the means of any two distributions is at least $\Omega(k)$ standard deviations. These results generally make heavy use of the generative model and particular properties of the distributions. In this paper, we show that a simple clustering algorithm works without assuming any generative (probabilistic) model. Our only assumption is what we call a “proximity condition”: the projection of any data point onto the line joining its cluster center to any other cluster center is $\Omega(k)$ standard deviations closer to its own center than the other center. Here the notion of standard deviations is based on the spectral norm of the matrix whose rows represent the difference between a point and the mean of the cluster to which it belongs. We show that in the generative models studied, our proximity condition is satisfied and so we are able to derive most known results for generative models as corollaries of our main result. We also prove some new results for generative models - e.g., we can cluster all but a small fraction of points only assuming a bound on the variance. Our algorithm relies on the well known k -means algorithm, and along the way, we prove a result of independent interest – that the k -means algorithm converges to the “true centers” even in the presence of spurious points provided the initial (estimated) centers are close enough to the corresponding actual centers and all but a small fraction of the points satisfy the proximity condition. Finally, we present a new technique for boosting the ratio of inter-center separation to standard deviation. This allows us to prove results for learning mixture of a class of distributions under weaker separation conditions.

*This work was done while the author was visiting Microsoft Research India Lab.

1 Introduction

Clustering is in general a hard problem. But, there has been a lot of research (see Section 3 for references) on proving that if we have data points generated by a mixture of k probability distributions, then one can cluster the data points into the k clusters, one corresponding to each component, provided the means of the different components are well-separated. There are different notions of well-separated, but mainly, the (best known) results can be qualitatively stated as:

“If the means of every pair of densities are at least $\text{poly}(k)$ times standard deviations apart, then we can learn the mixture in polynomial time.”

These results generally make heavy use of the generative model and particular properties of the distributions (Indeed, many of them specialize to Gaussians or independent Bernoulli trials). In this paper, we make no assumptions on the generative model of the data. We are still able to derive essentially the same result (loosely stated for now as):

“If the projection of any data point onto the line joining its cluster center to any other cluster center is $\Omega(k)$ times standard deviations closer to its own center than the other center (we call this the “proximity condition”), then we can cluster correctly in polynomial time.”

First, if the n points to be clustered form the rows of an $n \times d$ matrix A and C is the corresponding matrix of cluster centers (so each row of C is one of k vectors, namely the centers of k clusters) then note that the maximum directional variance (no probabilities here, the variance is just the average squared distance from the center) of the data in any direction is just

$$\frac{1}{n} \cdot \text{Max}_{v:|v|=1} |(A - C) \cdot v|^2 = \frac{\|A - C\|^2}{n},$$

where $\|A - C\|$ is the spectral norm. So, spectral norm scaled by $1/\sqrt{n}$ will play the role of standard deviation in the above assertion. To our knowledge, this is the first result proving that clustering can be done in polynomial time in a general situation with only deterministic assumptions. It settles an open question raised in [KV09].

We will show that in the generative models studied, our proximity condition is satisfied and so we are able to derive all known results for generative models as corollaries of our theorem (with one qualification: whereas our separation is in terms of the whole data variance, often, in the case of Gaussians, one can make do with separations depending only on individual densities’ variances – see Section 3.)

Besides Gaussians, the planted partition model (defined later) has also been studied; both these distributions have very “thin tails” and a lot of independence, so one can appeal to concentration results. In section 6.3, we give a clustering algorithm for a mixture of general densities for which we only assume bounds on the variance (and no further concentration). Based on our algorithm, we show how to classify all but an ε fraction of points in this model. Section 3 has references to recent work dealing with distributions which may not even have variance, but these results are only for the special class of product densities, with additional constraints.

One crucial technical result we prove (Theorem 5.5) may be of independent interest. It shows that the good old k -means algorithm [Llo82] converges to the “true centers” even in the presence of spurious points provided the initial (estimated) centers are close enough to the corresponding actual centers and all but an ε fraction of the points satisfy the proximity condition. Convergence (or lack of it) of the k -means algorithm is again well-studied ([ORSS06, AV06, Das03, HPS05]). The result of [ORSS06] (one of the few to formulate sufficient conditions for the k -means algorithm to provably work) assumes the condition that the optimal clustering with k centers is substantially better than that with fewer centers and shows that one iteration of k -means yields a near-optimal solution. We show in section 6.4 that their condition implies

proximity for all but an ε fraction of the points. This allows us to prove that our algorithm, which is again based on the k -means algorithm, gives a PTAS.

The proof of Theorem 5.5 is based on Theorem 5.4 which shows that if current centers are close to the true centers, then misclassified points (whose nearest current center is not the one closest to the true center) are far away from true centers and so there cannot be too many of them. This is based on a clean geometric argument shown pictorially in Figure 2. Our main theorem in addition allows for an ε fraction of “spurious” points which do not satisfy the proximity condition. Such errors have often proved difficult to account for.

As indicated, all results on generative models assume a lower bound on the inter-center separation in terms of the spectral norm. In section 7, we describe a construction (when data is from a generative model – a mixture of distributions) which boosts the ratio of inter-center separation to spectral norm. The construction is the following: we pick two sets of samples A_1, A_2, \dots, A_n and B_1, B_2, \dots, B_n independently from the mixture. We define new points X_1, X_2, \dots, X_n , where X_i is defined as $(A'_i \cdot B'_1, A'_i \cdot B'_2, \dots, A'_i \cdot B'_n)$, where $'$ denotes that we have subtracted the mean (of the mixture.) Using this, we are able to reduce the dependence of inter-center separation on the minimum weight of a component in the mixture that all models generally need. This technique of boosting is likely to have other applications.

2 Preliminaries and the Main Theorem

For a matrix A , we shall use $\|A\|$ to denote its spectral norm. For a vector v , we use $|v|$ to denote its length. We are given n points in \mathbb{R}^d which are divided into k clusters – T_1, T_2, \dots, T_k . Let μ_r denote the mean of cluster T_r and n_r denote $|T_r|$. Let A be the $n \times d$ matrix with rows corresponding to the points. Let C be the $n \times d$ matrix where $C_i = \mu_r$, for all $i \in T_r$. We shall use A_i to denote the i^{th} row of A . Let

$$\Delta_{rs} = \left(\frac{ck}{\sqrt{n_r}} + \frac{ck}{\sqrt{n_s}} \right) \|A - C\|,$$

where c is a large enough constant.

Definition 2.1 *We say a point $A_i \in T_r$ satisfies the proximity condition if for any $s \neq r$, the projection of A_i onto the μ_r to μ_s line is at least Δ_{rs} closer to μ_r than to μ_s . We let G (for good) be the set of points satisfying the proximity condition.*

Note that the proximity condition implies that the distance between μ_r and μ_s must be at least Δ_{rs} . We are now ready to state the theorem.

Theorem 2.2 *If $|G| \geq (1 - \varepsilon) \cdot n$, then we can correctly classify all but $O(k^2 \varepsilon \cdot n)$ points in polynomial time. In particular, if $\varepsilon = 0$, all points are classified correctly.*

Often, when applying this theorem to learning a mixture of distributions, A will correspond to a set of n independent samples from the mixture. We will denote the corresponding distributions by F_1, \dots, F_k , and their relative weights by w_1, \dots, w_k . Often, σ_r will denote the maximum variance along any direction of the distribution F_r , and σ will denote $\max_r \sigma_r$. We denote the minimum mixing weight of a distribution as w_{\min} .

3 Previous Work

Learning mixture of distributions is one of the central problems in machine learning. There is vast amount of literature on learning mixture of Gaussian distributions. One of the most popular methods for this is the well

known EM algorithm which maximizes the log likelihood function [DLR77]. However, there are few results which demonstrate that it converges to the optima solution. Dasgupta [Das99] introduced the problem of learning distributions under suitable *separation conditions*, i.e., we assume that the distance between the means of the distributions in the mixture is large, and the goal is to recover the original clustering of points (perhaps with some error).

We first summarize known results for learning mixtures of Gaussian distributions under separation conditions. We ignore logarithmic factors in separation condition. We also ignore the minimum number of samples required by the various algorithms – they are often bounded by a polynomial in the dimension and the mixing weights. Let σ_r be the maximum variance of the Gaussian F_r in any direction. Dasgupta [Das99] gave an algorithm based on random projection to learn mixture of Gaussians provided mixing weights of all distributions are about the same, and $|\mu_i - \mu_j|$ is $\Omega((\sigma_i + \sigma_j) \cdot \sqrt{n})$. Dasgupta and Schulman [DS07] gave an EM based algorithm provided $|\mu_i - \mu_j|$ is $\Omega((\sigma_i + \sigma_j) \cdot n^{\frac{1}{4}})$. Arora and Kannan [AK01] also gave a learning algorithm with similar separation conditions. Vempala and Wang [VW04] were the first to demonstrate the effectiveness of spectral techniques. For spherical Gaussians, their algorithm worked with a much weaker separation condition of $\Omega((\sigma_i + \sigma_j) \cdot k^{\frac{1}{4}})$ between μ_i and μ_j . Achlioptas and McSherry [AM05] extended this to arbitrary Gaussians with separation between μ_i and μ_j being at least $\tilde{\Omega}\left(\left(k + \frac{1}{\sqrt{\min(w_i, w_j)}}\right) \cdot (\sigma_i + \sigma_j)\right)$. Kannan et. al. [KSV08] also gave an algorithm for arbitrary Gaussians with the corresponding separation being $\Omega\left(\frac{k^{\frac{3}{2}}}{w_{\min}^2} \cdot (\sigma_i + \sigma_j)\right)$. Recently, Brubaker and Vempala [BV08] gave a learning algorithm where the separation only depends on the variance perpendicular to a hyperplane separating two Gaussians (the so called “parallel pancakes problem”).

Much less is known about learning mixtures of heavy tailed distributions. Most of the known results assume that each distribution is a product distribution, i.e., projection along co-ordinate axes are independent. Often, they also assume some *slope condition* on the line joining any two means. These slope conditions typically say that the unit vector along such lines does not lie almost entirely along very few coordinates. Such a condition is necessary because if the only difference between two distributions were a single coordinate, then one would require much stronger separation conditions. Dasgupta et. al. [DHKS05] considered the problem of learning product distributions of heavy tailed distributions when each component distribution satisfied the following mild condition : $P[|X - \mu| \geq \alpha R] \leq \frac{1}{2\alpha}$. Here R is the half-radius of the distribution (these distributions can have unbounded variance). Their algorithm could classify at least $(1 - \varepsilon)$ fraction of the points provided the distance between any two means is at least $\Omega\left(\frac{\sigma \cdot k^{\frac{5}{2}}}{\varepsilon^2}\right)$. Here R is the maximum half-radius of the distributions along any coordinate. Under even milder assumptions on the distributions and a slope condition, they could correctly classify all but ε fraction of the points provided the corresponding separation was $\Omega\left(\sigma \cdot \sqrt{\frac{k}{\varepsilon}}\right)$. Their algorithm, however, requires exponential (in d and k) amount of time. This problem was resolved by Chaudhuri and Rao [CR08]. Dasgupta et. al. [DHKM07] considered the problem of classifying samples from a mixture of arbitrary distributions with bounded variance in any direction. They showed that if the separation between the means is $\Omega(\sigma k)$ and a suitable slope condition holds, then all the samples can be correctly classified. Their paper also gives a general method for bounding the spectral norm of a matrix when the rows are independent (and some additional conditions hold). We will mention this condition formally in Section 6 and make heavy use of it.

Finally, we discuss the *planted partition model* [McS01]. In this model, an instance consists of a set of n points, and there is an implicit partition of these n points into k groups. Further, there is an (unknown) $k \times k$ matrix of probabilities P . We are given a graph G on these n points, where an edge between two vertices from groups i and j is present with probability P_{ij} . The goal is to recover the actual partition of the points (and hence, an approximation to the matrix P as well). We can think of this as a special case of

learning mixture of k distributions, where the distribution F_r corresponding to the r^{th} part is as follows : F_r is a distribution over $\{0, 1\}^n$, one coordinate corresponding to each vertex. The coordinate corresponding to vertex u is set to 1 with probability P_{ij} , where j denotes the group to which u belongs. Note that the mean of F_r , μ_r , is equal to the vector $(P_{r\psi(u)})_{u \in V}$, where $\psi(u)$ denotes the group to which the vertex u belongs. McSherry [McS01] showed that if the following separation condition is satisfied, then one can recover the actual partition of the vertex set with probability at least $1 - \delta$ – for all $r, s, r \neq s$

$$|\mu_r - \mu_s|^2 \geq c \cdot \sigma^2 \cdot k \cdot \left(\frac{1}{w_{\min}} + \log \frac{n}{\delta} \right), \quad (1)$$

where c is a large constant, w_{\min} is such that every group has size at least $w_{\min} \cdot n$, and σ^2 denotes $\max_{i,j} P_{ij}$.

There is a rich body of work on the k -means problem and heuristic algorithms for this problem (see for example [KSS10, ORSS06] and references therein). One of the most widely used algorithms for this problem was given by Lloyd [Llo82]. In this algorithm, we start with an arbitrary set of k candidate centers. Each point is assigned to the closest candidate center – this clusters the points into k clusters. For each cluster, we update the candidate center to the mean of the points in the cluster. This gives a new set of k candidate centers. This process is repeated till we get a local optimum. This algorithm may take superpolynomial time to converge [AV06]. However, there is a growing body of work on proving that this algorithm gives a good clustering in polynomial time if the initial choice of centers is good [AV07, ADK09, ORSS06]. Ostrovsky et. al. [ORSS06] showed that a modification of the Lloyd’s algorithm gives a PTAS for the k -means problem if there is a sufficiently large separation between the means. Our result also fits in this general theme – the k -means algorithm on a choice of centers obtained from a simple spectral algorithm classifies the point correctly.

4 Our Contributions

Our main contribution is to show that a set of points satisfying a deterministic proximity condition (based on spectral norm) can be correctly classified (Theorem 2.2). The algorithm is described in Figure 1. It has two main steps – first find an initial set of centers based on SVD, and then run the standard k -means algorithm with these initial centers as seeds. In Section 5, we show that after each iteration of the k -means algorithm, the set of centers come exponentially close to the true centers. Although both steps of our algorithm – SVD and the k -means algorithm – have been well studied, ours is the first result which shows that *combining* the two leads to a provably good algorithm. In Section 6, we give several applications of Theorem 2.2. We have the following results for learning mixture of distributions (we ignore poly-logarithmic factors in the discussion below) :

- Arbitrary Gaussian Distributions with separation $\Omega\left(\frac{\sigma k}{\sqrt{w_{\min}}}\right)$: as mentioned above, this matches known results [AM05, KSV08] except for the fact that the separation condition between two distributions depends on the maximum standard deviation (as compared to standard deviations of these distributions only).
- Planted distribution model with separation $\Omega\left(\frac{k\sigma}{\sqrt{w_{\min}}}\right)$: this matches the result of McSherry [McS01] except for a \sqrt{k} factor which we can also remove with a more careful analysis.
- Distributions with bounded variance along any direction : we can classify all but an ε fraction of points if the separation between means is at least $\Omega\left(\frac{k\sigma}{\sqrt{\varepsilon}}\right)$. Although results are known for classifying (all but a small fraction) points from mixtures of distributions with unbounded variance [DHKS05, CR08], such results work for product distributions only.

- PTAS using the k -means algorithm : We show that the separation condition of Ostrovsky et. al. [ORSS06] is stronger than the proximity condition. Using this fact, we are also able to give a PTAS based on the k -means algorithm.

Further, ours is the first algorithm which applies to all of the above settings. In Section 7, we give a general technique for working with weaker separation conditions (for learning mixture of distributions). Under certain technical conditions described in Section 7, we give a construction which increases the spectral norm of $A - C$ at a much faster rate than the increase in inter-mean distance as we increase the number of samples. As applications of this technique, we have the following results :

- Arbitrary Gaussians with separation $\Omega\left(\sigma k \cdot \log \frac{d}{w_{\min}}\right)$: this is the first result for arbitrary Gaussians where the separation depends only logarithmically on the minimum mixture weight.
- Power-law distributions with sufficiently large (but constant) exponent γ (defined in equation (13)) : We prove that we can learn all but ε fraction of samples provided the separation between means is $\Omega\left(\sigma k \cdot \left(\log \frac{d}{w_{\min}} + \frac{1}{\varepsilon^{\frac{1}{\gamma}}}\right)\right)$. For large values of γ , it significantly reduces the dependence on ε .

We expect this technique to have more applications.

5 Proof of Theorem 2.2

Our algorithm for correctly classifying the points will run in several iterations. At the beginning of each iteration, it will have a set of k candidate points. By a Lloyd like step, it will replace these points by another set of k points. This process will go on for polynomial number of steps.

1. **(Base case)** Let \hat{A}_i denote the projection of the points on the best k -dimensional subspace found by computing SVD of A . Let $\nu_r, r = 1, \dots, k$, denote the centers of a (near)-optimal solution to the k -means problem for the points \hat{A}_i .
2. For $\ell = 1, 2, \dots$ do
 - (i) Assign each point A_i to the closest point among $\nu_r, r = 1, \dots, k$. Let S_r denote the set of points assigned to ν_r .
 - (ii) Define η_r as the mean of the points S_r . Update $\eta_r, r = 1, \dots, k$ as the new centers, i.e., set $\nu_r = \eta_r$ for the next iteration.

Figure 1: Algorithm Cluster

The iterative procedure is described in Figure 1. In the first step, we can use any constant factor approximation algorithm for the k -means problem. Note that the algorithm is same as Lloyd's algorithm, but we start with a special set of initial points as described in the algorithm. We now prove that after the first step (the base case), the estimated centers are close to the actual ones – this case follows from [KV09], but we prove it below for sake of completeness.

Lemma 5.1 (Base Case) *After the first step of the algorithm above,*

$$|\mu_r - \nu_r| \leq 20\sqrt{k} \cdot \frac{\|A - C\|}{\sqrt{n_r}}.$$

Proof. Suppose, for sake of contradiction, that there exists an r such that all the centers ν_1, \dots, ν_k are at least $\frac{20\sqrt{k} \cdot \|A - C\|}{\sqrt{n_r}}$ distance away from μ_r . Consider the points in T_r . Suppose $A_i \in T_r$ is assigned to the center $\nu_{c(i)}$ in this solution. The assignment cost for these points in this optimal k -means solution is

$$\begin{aligned} \sum_{i \in T_r} |\hat{A}_i - \nu_{c(i)}|^2 &= \sum_{i \in T_r} |(\mu_r - \nu_{c(i)}) - (\mu_r - \hat{A}_i)|^2 \\ &\geq \frac{|T_r|}{2} \cdot \left(\frac{20\sqrt{k} \cdot \|A - C\|}{\sqrt{n_r}} \right)^2 - \sum_{i \in T_r} |\mu_r - \hat{A}_i|^2 \end{aligned} \quad (2)$$

$$\geq 20k \cdot \|A - C\|^2 - 5k \cdot \|A - C\|^2 = 15k \|A - C\|^2 \quad (3)$$

where inequality (2) follows from the fact that for any two numbers a, b , $(a-b)^2 \geq \frac{a^2}{2} - b^2$; and inequality (3) follows from the fact that $\|\hat{A} - C\|_F^2 \leq 5k \cdot \|A - C\|^2$. But this is a contradiction, because one feasible solution to the k -means problem is to assign points in $\hat{A}_i, i \in T_s$ to μ_s for $s = 1, \dots, k$ – the cost of this solution is $\|\hat{A} - C\|_F^2 \leq 5k \|A - C\|^2$. ■

Observe that the lemma above implies that there is a unique center ν_r associated with each μ_r . We now prove a useful lemma which states that removing small number of points from a cluster T_r can move the mean of the remaining points by only a small distance.

Lemma 5.2 *Let X be a subset of T_r . Let $m(X)$ denote the mean of the points in X . Then*

$$|m(X) - \mu_r| \leq \frac{\|A - C\|}{\sqrt{|X|}}.$$

Proof. Let u be unit vector along $m(X) - \mu_r$. Now,

$$|(A - C) \cdot u| \geq \left(\sum_{i \in X} ((A_i - \mu_r) \cdot u)^2 \right)^{\frac{1}{2}} \geq \frac{1}{\sqrt{|X|}} \left(\sum_{i \in X} |(A_i - \mu_r) \cdot u| \right) \geq \sqrt{|X|} \cdot |m(X) - \mu_r|$$

But, $|(A - C) \cdot u| \leq \|A - C\|$. This proves the lemma. ■

Corollary 5.3 *Let $Y \subseteq T_s$ such that $|T_s - Y| \leq \delta \cdot n_s$, where $\delta < \frac{1}{2}$. Let $m(Y)$ denote the mean of the points in Y . Then*

$$|m(Y) - \mu_s| \leq \frac{2 \cdot \sqrt{\delta} \cdot \|A - C\|}{\sqrt{n_s}}.$$

Proof. Let X denote $T_s - Y$. We know that $\mu_s \cdot |T_s| = |X| \cdot m(X) + |Y| \cdot m(Y)$. So we get

$$|m(Y) - \mu_s| = \frac{|X|}{|Y|} \cdot |m(X) - \mu_s| \leq \frac{\sqrt{|X|}}{|Y|} \cdot \|A - C\|$$

where the inequality above follows from Lemma 5.2. The result now follows because $|Y| \geq \frac{n_s}{2}$. ■

Now we show that if the estimated centers are close to the actual centers, then one iteration of the second step in the algorithm will reduce this separation by at least half.

Notation :

- $\nu_1, \nu_2, \dots, \nu_k$ denote the current centers at the beginning of an iteration in the second step of the algorithm, where ν_r is the current center closest to μ_r .

- S_r denotes the set of points A_i for which the closest current center is ν_r .
- η_r denotes the mean of points in S_r ; so η_r are the new centers. Let $\delta_r = |\mu_r - \nu_r|$.

The theorem below shows that the set of misclassified points (which really belong to T_r , but have $\nu_s, s \neq r$, as the closest current center) are not too many in number. The proof first shows that any misclassified point must be far away from μ_r and since the sum of squared distances from μ_r for all points in T_r is bounded, there cannot be too many.

Theorem 5.4 *Assume that $\delta_r + \delta_s \leq \Delta_{rs}/16$ for all $r \neq s$. Then,*

$$|T_r \cap S_s \cap G| \leq \frac{6ck \cdot \|A - C\|^2 (\delta_r^2 + \delta_s^2)}{\Delta_{rs}^2 |\mu_r - \mu_s|^2} \quad (4)$$

Further, for any $W \subseteq T_r \cap S_s$,

$$|m(W) - \mu_s| \leq \frac{100 \cdot \|A - C\|}{\sqrt{|W|}} \quad (5)$$

Proof. Let \bar{v} denote the projection of vector v to the affine space V spanned by $\mu_1, \dots, \mu_k, \nu_1, \dots, \nu_k$ and $\eta_1, \eta_2, \dots, \eta_k$. Assume $A_i \in T_r \cap S_s \cap G$. Splitting \bar{A}_i into its projection along the line μ_r to μ_s and the component orthogonal to it, we can write

$$\bar{A}_i = \frac{1}{2}(\mu_r + \mu_s) + \lambda(\mu_r - \mu_s) + u,$$

where u is orthogonal to $\mu_r - \mu_s$. Since \bar{A}_i is closer to ν_s than to ν_r , we have

$$\begin{aligned} \bar{A}_i \cdot (\nu_s - \nu_r) &\geq \frac{1}{2}(\nu_s - \nu_r) \cdot (\nu_r + \nu_s) \\ \text{i.e., } \frac{1}{2}(\mu_r + \mu_s) \cdot (\nu_s - \nu_r) + \lambda(\mu_r - \mu_s) \cdot (\nu_s - \nu_r) + u \cdot (\nu_s - \nu_r) &\geq \frac{1}{2}(\nu_s - \nu_r) \cdot (\nu_s + \nu_r). \end{aligned}$$

We have $u \cdot (\nu_s - \nu_r) = u \cdot ((\nu_s - \mu_s) - (\nu_r - \mu_r))$ since u is orthogonal to $\mu_r - \mu_s$. The last quantity is at most $|u|\delta$, where $\delta = \delta_r + \delta_s$. Substituting this we get

$$\begin{aligned} \frac{1}{2} \cdot (\mu_r + \mu_s - \nu_r - \nu_s) \cdot (\nu_s - \nu_r) + \lambda(\mu_r - \mu_s) \cdot (\nu_s - \nu_r) + |u| \cdot \delta &\geq 0 \\ \text{i.e., } \frac{\delta^2}{2} + \frac{\delta}{2}|\mu_r - \mu_s| - \lambda|\mu_r - \mu_s|^2 + \lambda\delta|\mu_r - \mu_s| + |u|\delta &\geq 0. \end{aligned} \quad (6)$$

Now,

$$\begin{aligned} |\bar{A}_i - \mu_r| &= \left| \left(\frac{1}{2} - \lambda \right) \cdot (\mu_s - \mu_r) + u \right| \\ &\geq |u| \geq \frac{\lambda}{\delta} \cdot |\mu_r - \mu_s|^2 - \lambda|\mu_r - \mu_s| - \frac{\delta}{2} - \frac{|\mu_r - \mu_s|}{2} \quad \text{using (6)} \\ &\geq \frac{\Delta_{rs}|\mu_r - \mu_s|}{64\delta}, \end{aligned}$$

where the last inequality follows from the fact that $\lambda \geq \frac{\Delta_{rs}}{2|\mu_r - \mu_s|}$ (proximity condition) and the assumption that $\delta \leq \Delta_{rs}/16$. Therefore, we have

$$|T_r \cap S_s \cap G| \cdot \frac{\Delta_{rs}^2 |\mu_r - \mu_s|^2}{c\delta^2} \leq \sum_{i \in T_r \cap S_s \cap G} |\bar{A}_i - \mu_r|^2 \leq \sum_{i \in T_r} |\bar{A}_i - C_i|^2.$$

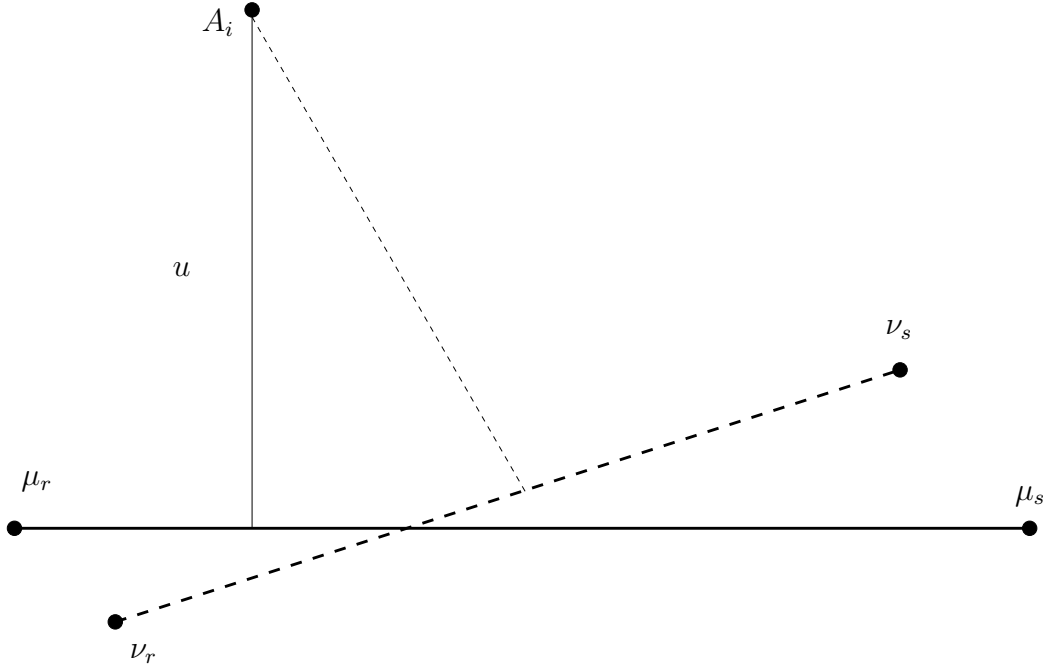


Figure 2: Misclassified A_i

If we take a basis u_1, u_2, \dots, u_p of V , we see that $\sum_{i \in T_r} |\bar{A}_i - C_i|^2 = \sum_{t=1}^p \sum_{i \in T_r} |(\bar{A}_i - C_i) \cdot u_t|^2 = \sum_{t=1}^p \|A - C\|^2 \leq 3k \|A - C\|^2$, which proves the first statement of the theorem.

For the second statement, we can write $m(W)$ as

$$m(W) = \frac{1}{2}(\mu_r + \mu_s) + \lambda(\mu_r - \mu_s) + u,$$

where, u is orthogonal to $\mu_r - \mu_s$. Since $m(W)$ is the average of points in S_s , we get (arguing as for (6)):

$$|u| \geq \frac{\lambda}{10\delta} |\mu_r - \mu_s|^2.$$

Now, we have

$$|m(W) - \mu_r|^2 = |u|^2 + \left(\lambda - \frac{1}{2}\right)^2 |\mu_r - \mu_s|^2, \text{ and } |m(W) - \mu_s|^2 = |u|^2 + \left(\lambda + \frac{1}{2}\right)^2 |\mu_r - \mu_s|^2$$

If $\lambda \leq 1/4$, then clearly, $|m(W) - \mu_s| \leq 4|m(W) - \mu_r|$. If $\lambda > 1/4$, then we have $|u| \geq \frac{\lambda}{10\delta} |\mu_r - \mu_s|^2 \geq \frac{1}{2} \cdot \left(\lambda + \frac{1}{2}\right) \cdot |\mu_r - \mu_s|$ because $\frac{|\mu_r - \mu_s|}{\delta} \geq 16$. This again yields $|m(W) - \mu_s| \leq 4|u| \leq 4|m(W) - \mu_r|$.

Now, by Lemma 5.2, we have $|m(W) - \mu_r| \leq \frac{\|A - C\|}{\sqrt{|W|}}$, so the second statement in the theorem. ■

We are now ready to prove the main theorem of this section which will directly imply Theorem 2.2. This shows that k -means converges if the starting centers are close enough to the corresponding true centers. To gain intuition, it is best to look at the case $\varepsilon = 0$, when all points satisfy the proximity condition. Then the theorem says that if $|\nu_s - \mu_s| \leq \frac{\gamma \|A - C\|}{\sqrt{n_s}}$ for all s , then $|\eta_s - \mu_s| \leq \frac{\gamma \|A - C\|}{2\sqrt{n_s}}$, thus halving the upper bound of the distance to μ_s in each iteration.

Theorem 5.5 *If*

$$\delta_s \leq \max \left(\frac{\gamma \cdot \|A - C\|}{\sqrt{n_s}}, 800\sqrt{k\varepsilon n} \cdot \frac{\|A - C\|}{n_s} \right),$$

for all s and a parameter $\gamma \leq ck/50$, then

$$|\eta_s - \mu_s| \leq \max \left(\frac{\gamma \cdot \|A - C\|}{2\sqrt{n_s}}, 800\sqrt{k\varepsilon n} \cdot \frac{\|A - C\|}{n_s} \right),$$

for all s .

Proof. Let n_{rs}, μ_{rs} denote the number and mean respectively of $T_r \cap S_s \cap G$ and n'_{rs}, μ'_{rs} of $(T_r \setminus G) \cap S_s$. Similarly, define n_{ss} and μ_{ss} as the size and mean of the points in $T_s \cap S_s$. We get

$$|S_s| \eta_s = n_{ss} \mu_{ss} + \sum_{r \neq s} n_{rs} \mu_{rs} + \sum_r n'_{rs} \mu'_{rs}.$$

We have

$$|\mu_{ss} - \mu_s| \leq \frac{\sqrt{|S_s| - n_{ss}}}{n_{ss}} \|A - C\|; |\mu_{rs} - \mu_s| \leq \frac{100\|A - C\|}{\sqrt{n_{rs}}}; |\mu'_{rs} - \mu_s| \leq \frac{100\|A - C\|}{\sqrt{n'_{rs}}},$$

where the first one is from Corollary 5.3 (it is easy to check from the first statement of Theorem 5.4 that $n_{ss} \geq n_s/2$) and the last two are from the second statement in Theorem 5.4.

Now using the fact that length is a convex function, we have

$$\begin{aligned} |\eta_s - \mu_s| &\leq \frac{n_{ss}}{|S_s|} |\mu_{ss} - \mu_s| + \sum_{r \neq s} \frac{n_{rs}}{|S_s|} |\mu_{rs} - \mu_s| + \sum_r \frac{n'_{rs}}{|S_s|} |\mu'_{rs} - \mu_s| \\ &\leq 200\|A - C\| \left(\frac{\sqrt{|S_s| - n_{ss}}}{n_s} + \sum_{r \neq s} \frac{\sqrt{n_{rs}}}{n_s} + \sum_r \frac{\sqrt{n'_{rs}}}{n_s} \right) \\ &\leq 400\|A - C\| \left(\sum_{r \neq s} \frac{\sqrt{n_{rs}}}{n_s} + \sum_r \frac{\sqrt{n'_{rs}}}{n_s} \right) \end{aligned}$$

since $|S_s| - n_{ss} = \sum_{r \neq s} n_{rs} + \sum_s n'_{rs}$. Let us look at each of the terms above. Note that $n_{rs} \leq \frac{24ck\|A-C\|^2 \max(\delta_r, \delta_s)^2}{\Delta_{rs}^2 |\mu_r - \mu_s|^2}$ (using Theorem 5.4). So

$$\begin{aligned} \sum_{r \neq s} \frac{\sqrt{n_{rs}}}{n_s} &\leq \frac{5\sqrt{ck}\|A - C\|}{n_s} \sum_{r \neq s} \frac{\max(\delta_r, \delta_s)}{\Delta_{rs} \cdot |\mu_r - \mu_s|} \\ &\leq \frac{5\sqrt{ck}\|A - C\|^2}{n_s} \sum_{r \neq s} \frac{1}{\Delta_{rs} \cdot |\mu_r - \mu_s|} \cdot \left(\frac{\gamma}{\sqrt{\min(n_r, n_s)}} + \frac{800\sqrt{k\varepsilon n}}{\min(n_r, n_s)} \right) \\ &\leq \frac{5\sqrt{ck}}{n_s} \sum_{r \neq s} \frac{\min(n_r, n_s)}{c^2 k^2} \cdot \left(\frac{\gamma}{\sqrt{\min(n_r, n_s)}} + \frac{800\sqrt{k\varepsilon n}}{\min(n_r, n_s)} \right) \\ &\leq \frac{\gamma}{c\sqrt{n_s}} + \frac{\sqrt{k\varepsilon n}}{n_s} \end{aligned}$$

Also, note that $\sum_r n'_{rs} \leq \varepsilon n$. So we get $\sum_r \frac{\sqrt{n'_{rs}}}{n_s} \leq \frac{\sqrt{k\varepsilon n}}{n_s}$. Assuming c to be large enough constant proves the theorem. ■

Now we can easily finish the proof of Theorem 2.2. Observe that after the base case in the algorithm, the statement of Theorem 5.5 holds with $\gamma = 20\sqrt{k}$. So after enough number of iterations of the second step in our algorithm, γ will become very small and so we will get

$$\delta_s \leq 800\sqrt{k\varepsilon n} \cdot \frac{\|A - C\|}{n_s},$$

for all s . Now substituting this in Theorem 5.4, we get

$$|T_r \cap S_s \cap G| \leq \frac{800^2 \cdot 24ck \cdot \|A - C\|^4 \cdot \varepsilon \cdot n \cdot k}{\Delta_{rs}^2 |\mu_r - \mu_s|^2 \cdot \min(n_r, n_s)^2} \leq \varepsilon n$$

Summing over all pairs r, s implies Theorem 2.2.

6 Applications

We now give applications of Theorem 2.2 to various settings. One of the main technical steps here would be to bound the spectral norm of a random $n \times d$ matrix Y whose rows are chosen independently. We use the following result from [DHKM07]. Let D denote the matrix $E[Y^T Y]$. Also assume that $n \geq d$.

Fact 6.1 *Let γ be such that $\max_i |Y_i| \leq \gamma\sqrt{n}$ and $\|D\| \leq \gamma^2 n$. Then $\|Y\| \leq \gamma \cdot \sqrt{n} \cdot \text{polylog}(n)$ with high probability.*

6.1 Learning in the planted distribution model

In this section, we show that McSherry's result [McS01] can be derived as a corollary of our main theorem. Consider an instance of the planted distribution model which satisfies the condition (1). We would like to show that with high probability, the points satisfy the proximity condition. Fix a point $A_i \in T_r$. We will show that the probability that it does not satisfy this condition is at most $\frac{\delta}{n}$. Using union bound, it will then follow that the proximity condition is satisfied with probability at least $1 - \delta$.

Let $s \neq r$. Let v denote the unit vector along $\mu_r - \mu_s$. Let L_{rs} denote the line joining μ_r and μ_s , and \hat{A}_i be the projection of A_i on L_{rs} . The following result shows that the distance between \hat{A}_i and μ_r is small with high probability.

Lemma 6.2 *Assume $\sigma \geq \frac{3 \log n}{n}$, where $\sigma = \max_{i,j} \sqrt{P_{ij}}$. With probability at least $1 - \frac{\delta}{n \cdot k}$,*

$$|\hat{A}_i - \mu_r| \leq ck \cdot \sigma \cdot \left(\log \left(\frac{n}{\delta} \right) + \frac{1}{\sqrt{w_{\min}}} \right),$$

where c is a large constant.

Proof. For a vector A_i , we use A_{ij} to denote the coordinate of A_i at position j . Define μ_{rj} similarly. First observe that $|\hat{A}_i - \mu_r| = |v \cdot (A_i - \mu_r)|$. The coordinates of v corresponding to points belonging to a particular cluster are same – let v^t denote this value for cluster T_t . So we get

$$|v \cdot (A_i - \mu_r)| \leq \sum_{t=1}^k |v^t| \cdot \left| \sum_{j \in T_t} (A_{ij} - \mu_{rj}) \right| \leq \sum_{t=1}^k \frac{\left| \sum_{j \in T_t} A_{ij} - P_{rt} \cdot n_t \right|}{\sqrt{n_t}}$$

where n_t denotes the size of cluster T_t . The last inequality above follows from the fact that $|v| = 1$. So, $1 \geq \sum_{j \in T_t} v_j^2 = (v^t)^2 \cdot n_t$. Now, if A_i does not satisfy the condition of the lemma, then there must be some t for which

$$\left| \sum_{j \in T_t} A_{ij} - P_{rt} \cdot n_t \right| \geq c\sigma\sqrt{n_t} \cdot \left(\log \left(\frac{n}{\delta} \right) + \frac{1}{\sqrt{w_{\min}}} \right)$$

Now note that $A_{ij}, j \in T_r$ are i.i.d. 0-1 random variables with mean P_{rt} . Now we use the following version of Chernoff bound : let X_1, \dots, X_l be i.i.d. 0-1 random variables, each with mean p . Then

$$\Pr \left[\left| \sum_{i=1}^l X_i - l \cdot p \right| \geq \eta \cdot lp \right] \leq \begin{cases} e^{-\eta^2 lp/4} & \text{if } \eta \leq 2e - 1 \\ 2^{-\eta \cdot lp} & \text{otherwise} \end{cases},$$

For us, $\eta = \frac{c\sigma}{P_{rt}\sqrt{n_t}} \cdot \left(\log \left(\frac{n}{\delta} \right) + \frac{1}{\sqrt{w_{\min}}} \right)$. If $\eta \leq 2e - 1$, the probability of this event is at most

$$\exp \left(-\frac{c^2\sigma^2}{n_t P_{rt}^2} \cdot n_t P_{rt} \cdot \log \left(\frac{n}{\delta} \right) \right) \leq \frac{\delta}{n^2}.$$

Now, assume $\eta > 2e - 1$. In this case the probability of this event is at most

$$2^{-c\sigma\sqrt{\frac{n_t}{w_{\min}}}} \leq \frac{1}{n^3},$$

where we have assumed that $\sigma \geq \frac{3 \log n}{n}$ (we need this assumption anyway to use Wigner's theorem for bounding $\|A - C\|$). ■

Assuming

$$|\mu_r - \mu_s| \geq 4ck \cdot \sigma \cdot \left(\log \left(\frac{n}{\delta} \right) + \frac{1}{\sqrt{w_{\min}}} \right),$$

we see that

$$|\hat{A}_i - \mu_s| - |\hat{A}_i - \mu_r| \geq \frac{ck\|A - C\|}{\sqrt{n_r}} + \frac{ck\|A - C\|}{\sqrt{n_s}},$$

with probability at least $1 - \frac{\delta}{nk}$. Here, we have used the fact that $\|A - C\| \leq c' \cdot \sigma\sqrt{n}$ with high probability (Wigner's theorem). Now, using union bound, we get that all the points satisfy the proximity condition with probability at least $1 - \delta$.

Remark : Here we have used C as the matrix whose rows are the actual means μ_r . But while applying Theorem 2.2, C should represent the means of the samples in A belonging to a particular cluster. The error incurred here can be made very small and will not affect the results. So we shall assume that μ_r is the actual mean of points in T_r . Similar comments apply in other applications described next.

6.2 Learning Mixture of Gaussians

We are given a mixture of k Gaussians F_1, \dots, F_k in d dimensions. Let the mixture weights of these distributions be w_1, \dots, w_k and μ_1, \dots, μ_k denote their means respectively.

Lemma 6.3 *Suppose we are given a set of $n = \text{poly} \left(\frac{d}{w_{\min}} \right)$ samples from the mixture distribution. Then these points satisfy the proximity condition with high probability if*

$$|\mu_r - \mu_s| \geq \frac{ck\sigma_{\max}}{\sqrt{w_{\min}}} \text{polylog} \left(\frac{d}{w_{\min}} \right),$$

for all $r, s, r \neq s$. Here σ_{\max} is the maximum variance in any direction of any of the distributions F_r .

Proof. It can be shown that $\|A - C\|$ is $O\left(\sigma_{\max}\sqrt{n} \cdot \text{polylog}\left(\frac{d}{w_{\min}}\right)\right)$ with high probability (see [DHKM07]). Further, let p be a point drawn from the distribution F_r . Let L_{rs} be the line joining μ_r and μ_s . Let \hat{p} be the projection of p on this line. Then the fact that F_r is Gaussian implies that $|\hat{p} - \mu_r| \leq \sigma_{\max}\text{polylog}(n)$ with probability at least $1 - \frac{1}{n^2}$. It is also easy to check that the number of points from F_r in the sample is close to $w_r n$ with high probability. Thus, it follows that all the points satisfy the proximity condition with high probability. ■

The above lemma and Theorem 2.2 imply that we can correctly classify all the points. Since we shall sample at least $\text{poly}(d)$ points from each distribution, we can learn each of the distribution to high accuracy.

6.3 Learning Mixture of Distributions with Bounded Variance

We consider a mixture of distributions F_1, \dots, F_k with weights w_1, \dots, w_k . Let σ be an upper bound on the variance along any direction of a point sampled from one of these distributions. In other words,

$$\sigma \geq E\left[\left((x - \mu_r) \cdot v\right)^2\right],$$

for all distributions F_r and each unit vector v .

Theorem 6.4 *Suppose we are given a set of $n = \text{poly}\left(\frac{d}{w_{\min}}\right)$ samples from the mixture distribution. Assume that $\sigma \geq \frac{\text{polylog}(n)}{\sqrt{d}}$. Then there is an algorithm to correctly classify at least $1 - \varepsilon$ fraction of the points provided*

$$|\mu_r - \mu_s| \geq \frac{40k\sigma}{\sqrt{\varepsilon}} \text{polylog}\left(\frac{d}{\varepsilon}\right),$$

for all $r, s, r \neq s$. Here ε is assumed to be less than w_{\min} .

Proof. The algorithm is described in Figure 3. We now prove that this algorithm has the desired properties. Let A denote the $n \times d$ matrix of points and C be the corresponding matrix of means. We first bound the spectral norm of $A - C$. The bound obtained is quite high, but is probably tight.

1. Run the first step of Algorithm `Cluster` on the set of points, and let ν_1, \dots, ν_k denote the centers obtained.
2. Remove centers ν_r (and points associated with them) to which less than $d^2 \log d$ points are assigned. Let $\nu_1, \dots, \nu_{k'}$ be the remaining centers.
3. Remove any point whose distance from the nearest center in $\nu_1, \dots, \nu_{k'}$ is more than $\frac{\sigma\sqrt{n}}{\sqrt{d}}$.
4. Run the algorithm `Cluster` on the remaining set of points and output the clustering obtained.

Figure 3: Algorithm for Clustering points from mixture of distributions with bounded variance.

Lemma 6.5 *With high probability, $\|A - C\| \leq \sigma\sqrt{dn} \cdot \text{polylog}(n)$.*

Proof. We use Fact 6.1. Let Y denote $A - C$. Note that $|Y_i|^2 = \sum_{j=1}^d (A_{ij} - C_{ij})^2$. Since $E[(A_{ij} - C_{ij})^2] \leq \sigma^2$ (because it is the variance of this distribution along one of the coordinate axes), the expected value of

$|Y_i|^2$ is at most $\sigma^2 d$. Now using Chebychev's inequality, we see that $\max_i |Y_i| \geq \sigma \sqrt{dn} \text{polylog}(n)$ is at most $\frac{1}{\text{polylog}(n)}$. Now we consider $Y^T Y = \sum_i E[Y_i^T Y_i]$ (recall that we are treating Y_i as a row vector). So if v is a unit (column) vector, then $v^T E[Y^T Y] v = \sum_i E[|Y_i \cdot v|^2]$. But $E[|Y_i \cdot v|^2]$ is just the variance of the distribution corresponding to A_i along v . So this quantity is at most σ^2 for all v . Thus, we see that $\|E[Y^T Y]\| \leq \sigma^2 n$. This proves the lemma. ■

The above lemma allows us to bound the distance between μ_r and the nearest mean obtained in Step 1 of the algorithm. The proof proceeds along the same lines as that of Lemma 5.1.

Lemma 6.6 *For each μ_r , there exists a center ν_r such that ν_r is not removed in Step 2 and $|\mu_r - \nu_r| \leq \frac{10\sigma\sqrt{dk}}{\sqrt{\varepsilon}} \cdot \text{polylog}(n)$.*

Proof. Suppose the statement of the lemma is false for μ_r . At most $kd^2 \log d$, which is much less than $|T_r|$, points are assigned to a center which is removed in Step 2. The remaining points in T_r are assigned to centers which are not removed. So, arguing as in the proof of Lemma 5.1, the clustering cost in step 1 for points in T_r is at least

$$\begin{aligned} \frac{|T_r| - kd^2 \log d}{2} \cdot \left(\frac{10\sigma\sqrt{dk}}{\sqrt{\varepsilon}} \cdot \text{polylog}(n) \right)^2 - \sum_{i \in T_r} (\mu_r - \hat{A}_i)^2 &\geq 50\sigma^2 dkn \cdot \text{polylog}(n) - \|A - C\|^2 \\ &\geq 10k\|A - C\|^2 \end{aligned}$$

where the last inequality follows from Lemma 6.5. But, as in the proof of Lemma 5.1, this is a contradiction. ■

Note that in the lemma above, ν_r may not be unique for different means μ_r . Call a point $A_i \in T_r$ *bad* if $|A_i - \mu_r| \geq \frac{\sigma\sqrt{n}}{2}$. Call a point $A_i \in T_r$ *nice* if $|A_i - \mu_r| \leq \frac{\sigma\sqrt{n}}{2\sqrt{d}}$.

Lemma 6.7 *The number of bad points is at most $d \cdot \log d$ with high probability. The number of points which are not nice is at most $d^2 \log d$ with high probability. The number of nice points that are removed is at most $4kd^2 \log d$.*

Proof. Arguing as in the proof of Lemma 6.5, the probability that $|A_i - C_i| \geq \sigma \cdot \sqrt{n}$ is at most $\frac{d}{n}$. So the expected number of bad points is at most d . The first statement in the lemma now follows from Chernoff bound. The second statement is proved similarly. At most $kd^2 \log d$ points are removed in Step 1. Now suppose A_i is nice. Then Lemma 6.6 implies that it will not be removed in Step 3 (using Lemma 6.5 and the fact that n is large enough). ■

Corollary 6.8 *With high probability the following event happens : suppose ν_r does not get removed in Step 2. Then there is a mean μ_r such that $|\mu_r - \nu_r| \leq \frac{2\sigma\sqrt{n}}{\sqrt{d}}$.*

Proof. Since ν_r is not removed, it has at least one nice point A_i assigned to it (otherwise it will have at most $d^2 \log d$ points assigned to it and it will be removed). The distance of A_i to the nearest mean μ_s is at most $\frac{\sigma\sqrt{n}}{2\sqrt{d}}$, and Lemma 6.6 implies that there is a center ν_s which is not removed and for which $|\mu_s - \nu_s| \leq \frac{\sigma\sqrt{n}}{2\sqrt{d}}$. So, $|\nu_s - A_i| \leq \frac{\sigma\sqrt{n}}{\sqrt{d}}$. Since ν_r is the closest center to A_i , $|\nu_r - A_i| \leq \frac{\sigma\sqrt{n}}{\sqrt{d}}$ as well. Now, $|\nu_r - \mu_s| \leq |\nu_r - A_i| + |A_i - \mu_s|$ and the result follows. ■

Let A' be the set of points which are remaining after the third step of our algorithm. We now define a new clustering T'_1, \dots, T'_k of points in A' . This clustering will be very close to the actual clustering T_1, \dots, T_k

and so it will be enough to correctly cluster a large fraction of the points according to this new clustering. We define

$$T'_r = \{A_i \in T_r : A_i \text{ is not bad and does not get removed}\} \cup \{A_i : A_i \text{ is a bad point which does not get removed and the nearest center among the actual centers is } \mu_r\}.$$

Let μ'_r be the mean of T'_r and C' be the corresponding matrix of means.

Lemma 6.9 *With high probability, for all r , $\|A' - C'\| \leq O(\sigma \cdot \sqrt{n} \cdot \text{polylog}(n))$, and $|\mu_r - \mu_{r'}| \leq \frac{10\sigma kd^2 \log d}{\varepsilon \sqrt{n}}$.*

Proof. We first prove the second statement. The points in T'_r contain all the points in T_r except for at most $5kd^2 \log d$ points (Lemma 6.7). First consider the points in $T_r - T'_r$ which are not bad. Since all these points are at distance less than $\sigma\sqrt{n}$ from μ_r , the removal of these points shifts the mean by at most $\frac{5\sigma k \sqrt{n} d^2 \log d}{|T_r|} \leq \frac{\sigma kd^2 \log d}{\varepsilon \sqrt{n}}$. Now T'_r may contain some bad points as well. First observe that any such bad point must be at most $\frac{3\sigma\sqrt{n}}{\sqrt{d}}$ away from μ_r . Indeed, the reason why we retained this bad point in Steps 2 and 3 is because it is at distance at most $\frac{\sigma\sqrt{n}}{\sqrt{d}}$ from ν_r from some r . So combined with Corollary 6.8, this statement is true. So these bad points can again shift the mean by a similar amount. This proves the second part of the lemma.

Now we prove the first part of the lemma. Break $A' - C'$ into two parts $-A'_B - C'_B$ and $A'_G - C'_G$ – the rows of $A' - C'$ which are from bad points and the remaining rows (good) respectively. $A'_B - C'_B$ has at most $d \log d$ rows, and each row (as argued above) has length at most $\frac{4\sigma\sqrt{n}}{\sqrt{d}}$. So $\|A'_B - C'_B\| \leq 3\sigma\sqrt{n} \log(d)$. Now consider $A'_G - C'_G$. Let C_G be the rows of the original matrix C corresponding to A'_G . Then $\|A'_G - C'_G\| \leq \|A'_G - C_G\| + \|C'_G - C_G\|$. Each row of $C_G - C'_G$ has length at most $\frac{10\sigma kd^2 \log d}{\varepsilon \sqrt{n}}$ and so its spectral norm is at most $\frac{10\sigma kd^2 \log d}{\varepsilon}$, which is quite small (doesn't involve n at all). So it remains to bound $\|A'_G - C_G\|$. Let Z be the rows of $A - C$ which correspond to points which are not bad. Note that the rows of Z are independent and have length at most $\sigma\sqrt{n}$. So applying Fact 6.1 and arguing as in Lemma 6.5, we can show that $\|Z\|$ is at most $\sigma\sqrt{n} \cdot \text{polylog}(n)$. Now observe that $A'_G - C_G$ is obtained by picking some rows of Z (by a random process), and so its spectral norm is at most that of $\|Z\|$. This proves the lemma. ■

We are now ready to prove the main theorem. We would like to recover the clustering C' (since C and C' agree on all but the bad points). We argue that at least $(1 - \varepsilon)$ fraction of the points satisfy the proximity condition. Indeed, it is easy to check that at least $(1 - \varepsilon)$ fraction of the points A_i are at distance at most $\frac{4\sigma\sqrt{d}}{\sqrt{\varepsilon}}$ from the corresponding mean μ_r and satisfy the proximity condition. Since the distance between μ_s and μ'_s is *very small* (dependent inversely on n), and A_i is only $\frac{4\sigma\sqrt{d}}{\sqrt{\varepsilon}}$ far from μ_r , it will satisfy the proximity condition for A', C' as well (provided n is large enough). ■

6.4 Sufficient conditions for convergence of k -means

As mentioned in Section 3, Ostrovsky et. al. [ORSS06] provided the first sufficient conditions under which they prove effectiveness of (a suitable variant of) the k -means algorithm. Here, we show that their conditions are (much) stronger than the proximity condition. We first describe their conditions. Let Δ_k be the optimal cost of the k -means problem (i.e., sum of distance squared of each point to nearest center) with k centers. They require:

$$\Delta_k \leq \varepsilon \Delta_{k-1}.$$

Claim 6.10 *The condition above implies the proximity condition for all but ε fraction of the points.*

Proof. Suppose the above condition is true. One way of getting a solution with $k - 1$ centers is to remove a center μ_r and move all points in T_r to the nearest other center μ_s . Now, their condition implies

$$|\mu_r - \mu_s|^2 \geq \frac{1}{n_r \cdot \varepsilon} \|A - C\|_F^2 \quad \forall s \neq r.$$

If some ε fraction of T_r do not satisfy the proximity condition, then the distance squared of each such point to μ_r is at least the distance squared along the line μ_r to μ_s which is at least $(1/4)|\mu_r - \mu_s|^2$ which is at least $\Omega(\|A - C\|_F^2/n_r\varepsilon)$. So even the assignment cost of such points exceeds $\|A - C\|_F^2$, the total cost, a contradiction. This proves the claim. ■

We now show that our algorithm gives a PTAS for the k -means problem.

Getting a PTAS Let T_1, \dots, T_k be the optimal clustering and μ_1, \dots, μ_k be the corresponding means. As before, n_r denotes the size of T_r . Let G be the set of points which satisfy the proximity condition (the good points). The above claim shows that $|G| \geq (1 - \varepsilon)n$. For simplicity, assume that *exactly* ε fraction of the points do not satisfy the proximity condition.

Let S_1, \dots, S_r be the clustering output by our algorithm. Let μ'_r be the mean of S_r . First observe that Theorem 5.5 implies that when our algorithm stops,

$$|\mu_r - \mu'_r| \leq \frac{c \cdot \sqrt{k\varepsilon n}}{n_r} \cdot \|A - C\|. \quad (7)$$

for some constant c . For a point A_i , let $\alpha(A_i)$ the square of its distance to the closest mean among μ_1, \dots, μ_k . Define $\alpha'(A_i)$ for the solution output by our algorithm similarly.

Claim 6.11 *If $A_i \notin G$, then $\alpha'(A_i) \leq (1 + O(\varepsilon)) \cdot \alpha(A_i)$.*

Proof. Suppose $A_i \in T_r$, and it does not satisfy the proximity condition for the pair μ_r, μ_s . Then it is easy to see that $\alpha(A_i) \geq (|\mu_r - \mu_s|/2 - \Delta_{rs})^2 \geq \frac{|\mu_r - \mu_s|^2}{4} \geq \frac{\|A - C\|_F^2}{4n_r\varepsilon}$. Let μ_t be the closest mean to A_i . Then

$$\begin{aligned} \alpha'(A_i) &\leq |A_i - \mu'_t|^2 \leq (1 + \varepsilon)|A_i - \mu_t|^2 + \left(\frac{1}{\varepsilon} + 1\right) \cdot |\mu_t - \mu'_t|^2 \\ &\leq (1 + \varepsilon)\alpha(A_i) + \frac{c^2 k n}{n_r^2} \cdot \|A - C\|^2 \leq (1 + O(\varepsilon))\alpha(A_i) \end{aligned}$$

where the second last inequality follows from equation (7). Note that the constant in $O(\varepsilon)$ above contains terms involving k and w_{\min} . ■

Claim 6.12 *If $A_i \in G$, but is mis-classified by our algorithm, then $\alpha'(A_i) \leq (1 + O(\varepsilon)) \cdot \alpha(A_i)$.*

Proof. Suppose $A_i \in G \cap T_r \cap S_s$. We use the machinery developed in the proof of Theorem 5.4. Define λ, u as in the proof of this theorem. Clearly, $\alpha'(A_i) \leq \alpha(A_i) + 2|\mu_r - \mu_s|^2$ (here we have also used equation (7)). But note that $\alpha(A_i) \geq |u|^2$. Now, $|u| \geq \frac{\Delta_{rs}|\mu_r - \mu_s|}{64\delta} = \Omega\left(\frac{|\mu_r - \mu_s|}{\sqrt{\varepsilon}}\right)$ (again using equation (7)). This implies the result. ■

Claim 6.13 *For all r ,*

$$\sum_{A_i \in G \cap T_r \cap S_r} \alpha'(A_i) \leq (1 + O(\sqrt{\varepsilon})) \cdot \sum_{A_i \in G \cap T_r \cap S_r} \alpha(A_i) + O(\sqrt{\varepsilon}) \cdot \sum_{A_i \notin G} \alpha(A_i).$$

Proof. Clearly,

$$|A_i - \mu'_r|^2 \leq (1 + \sqrt{\varepsilon} \cdot \beta) |A_i - \mu_r|^2 + \left(1 + \frac{1}{\sqrt{\varepsilon} \cdot \beta}\right) |\mu_r - \mu'_r|^2,$$

where β is a large constant in terms of $k, \frac{1}{w_{\min}}, c$. Summing this over all points in $T_r \cap S_r$, we get

$$\sum_{A_i \in G \cap T_r \cap S_r} \alpha'(A_i) \leq (1 + O(\sqrt{\varepsilon})) \sum_{A_i \in G \cap T_r \cap S_r} \alpha(A_i) + \frac{\sqrt{\varepsilon} \|A - C\|^2}{4k}$$

where the last inequality follows from (7) assuming β is large enough. But now, the proof of Claim 6.11, implies that $\sum_{A_i \notin G} \alpha(A_i) \geq \frac{\|A - C\|^2}{4}$. So we are done. ■

Now summing over all r in Claim 6.13 and using Claims 6.11, 6.12 implies that our algorithm is also a PTAS.

7 Boosting

Recall that the proximity condition requires that the distance between the means be polynomially dependent on $\frac{1}{w_{\min}}$ – this could be quite poor when one of the clusters is considerably smaller than the others. In this section, we try to overcome this obstacle for a special class of distributions.

Let F_1, \dots, F_k be a mixture of distributions in d dimensions. Let A be the $n \times d$ matrix of samples from the distribution and C be the corresponding matrix of centers. Let D_{\min} denote $\min_{r,s,r \neq s} |\mu_r - \mu_s|$. Then the key property that we desire from the mixture of distributions is as follows. The following conditions should be satisfied with high probability :

1. For all $r, s, r \neq s$,

$$|\mu_r - \mu_s| \geq \frac{10k \|A - C\|}{\sqrt{n}} \quad (8)$$

2. For all i ,

$$|A_i - C_i| \leq D_{\min} \cdot \sqrt{dn}^\alpha \text{polylog}(n), \quad (9)$$

where α is a small enough constant (something like 0.1 will suffice).

3. For all $r, s, r \neq s$,

$$\sum_{i \in T_r} [(A_i - \mu_r) \cdot v]^2 \leq \frac{|\mu_r - \mu_s|^2}{16} \cdot |T_r| \quad (10)$$

where v is the unit vector joining μ_r and μ_s . This condition is essentially saying that the average variance of points in T_r along v is bounded by $\frac{|\mu_r - \mu_s|}{\sqrt{4}}$.

The number of samples n will be a polynomial in $\frac{d}{w_{\min}}$. Recall that D_{\min} denotes $\min_{r,s,r \neq s} |\mu_r - \mu_s|$. To simplify the presentation, we assume that $|\mu_r - \mu_s| \leq D_{\min} \cdot \left(\frac{d}{w_{\min}}\right)^\beta$ for a constant β for all pairs r, s . We will later show how to get rid of this assumption. We now sample two sets of n points from this distribution – call these A and B . Assume that both A and B satisfy the conditions (9) and (10). For all r , we assume

that the mean of $A_i, i \in T_r$ is μ_r and $T_r \cap A$ has size $w_r \cdot n$. We assume the same for the points in B . The error caused by removing this assumption will not change our results. Let μ denote the overall mean of the points in A (or B). Note that $\mu = \sum_r w_r \mu_r$. We translate the points so that the overall mean is 0. In other words, define a translation f as $f(x) = x - \mu$. Let A'_i denote $f(A_i)$. Define B'_i similarly. We now define a set X of n points in n dimensions. The point X_i is defined as

$$(A'_i \cdot B'_1, \dots, A'_i \cdot B'_n).$$

The correspondence between X_i and A_i naturally defines a partitioning of X into k clusters. Let $S_r, r = 1, \dots, k$, denote these clusters. The mean θ_r of S_r is

$$(C'_r \cdot B'_1, \dots, C'_r \cdot B'_n),$$

where $C'_r = C_r - \mu$. Let Z denote the matrix of means of X , i.e., $Z_i = \theta_r$ if $X_i \in S_r$. We now show that this process *amplifies* the distance between the means θ_r by a much bigger factor than $\frac{\|X-Z\|}{\sqrt{n}}$.

Lemma 7.1 *For all $r, s, r \neq s$,*

$$|\theta_r - \theta_s| \geq \frac{|\mu_r - \mu_s|^2}{4} \cdot \sqrt{w_{\min}} \sqrt{n}.$$

Proof. First observe that $(\mu_r - \mu_s) \cdot (\mu_r - \mu_s) = (\mu_r - \mu) \cdot (\mu_r - \mu_s) - (\mu_s - \mu) \cdot (\mu_r - \mu_s)$. So at least one of $|(\mu_r - \mu) \cdot (\mu_r - \mu_s)|, |(\mu_s - \mu) \cdot (\mu_r - \mu_s)|$ must be at least $\frac{|\mu_r - \mu_s|^2}{2}$. Assume without loss of generality that this is so for $|(\mu_r - \mu) \cdot (\mu_r - \mu_s)|$. Now consider the coordinates i of $\theta_r - \theta_s$ corresponding to S_r . Such a coordinate will have value $(\mu_r - \mu_s) \cdot B'_i$. Therefore,

$$\begin{aligned} |\theta_r - \theta_s|^2 &\geq \sum_{i \in S_r} [(\mu_r - \mu_s) \cdot (B_i - \mu)]^2 \\ &\geq \frac{|S_r|}{2} [(\mu_r - \mu_s) \cdot (\mu - \mu_r)]^2 - \sum_{i \in S_r} [(\mu_r - \mu_s) \cdot (B_i - \mu_r)]^2 \\ &\geq \frac{|S_r|}{16} \cdot |\mu_r - \mu_s|^4 \end{aligned}$$

where the last inequality follows from (10). This proves the lemma. ■

Now we bound $\|X - Z\|$.

Lemma 7.2 *With high probability,*

$$\frac{\|X - Z\|}{\sqrt{n}} \leq D_{\min}^2 \cdot d \cdot n^{2\alpha} \cdot \left(\frac{d}{w_{\min}}\right)^\beta \cdot \text{polylog}(n)$$

Proof. Let Y denote the matrix $X - Z$, and D denote the matrix $E[Y^T Y]$ where the expectation is over the choice of A and B . We shall use Fact 6.1 to bound $\|Y\|$. Thus, we just need to bound $\max_i |Y_i|$ and $\|D\|$. Let γ denote $D_{\min}^2 \cdot d \cdot n^{2\alpha} \cdot \left(\frac{d}{w_{\min}}\right)^\beta \cdot \text{polylog}(n)$.

Claim 7.3 *For all i ,*

$$|Y_i| \leq \sqrt{n} \cdot \gamma$$

Proof. Suppose $X_i \in S_r$. Then the j^{th} coordinate of Y is $(A_i - \mu_r) \cdot (B_j - \mu) = (A_i - \mu_r) \cdot ((B_j - \mu_{r'}) + (\mu - \mu_{r'}))$, where r' is such that $B_j \in T_{r'}$. Now, condition (9) implies that $|A_i - \mu_r|, |B_j - \mu_{r'}| \leq D_{\min} \cdot \sqrt{dn}^\alpha \text{polylog}(n)$. This implies the claim. ■

Claim 7.4

$$\|D\| \leq \gamma^2 n$$

Proof. We can write $Y^T Y$ as $\sum_i Y_i^T Y_i$. Let v be any unit vector. Then $v^T E[Y^T Y] v = \sum_i E|Y_i \cdot v|^2$. For a fixed i , where $A_i \in T_r$,

$$\begin{aligned} E|Y_i \cdot v|^2 &= E \left(\sum_j v_j [(X_i - \mu_r) \cdot (Y_j - \mu)] \right)^2 \\ &= \sum_j v_j^2 E [(X_i - \mu_r) \cdot (Y_j - \mu)]^2 \end{aligned}$$

where the last inequality follows from the fact that expectation of Y_j is μ and $Y_j, Y_{j'}$ are independent if $j \neq j'$. Rest of the argument is as in Claim 7.3. ■

The above two claims combined with Fact 6.1 imply the lemma. ■

Now we pick n to be at least $\left(\frac{d}{w_{\min}}\right)^{8(\beta+1)}$. Assuming $\alpha < 0.1$, this implies (using the above two lemmas) that for all $r, s, r \neq s$,

$$|\theta_r - \theta_s| \geq \frac{\|X - Z\|}{\sqrt{n}} \cdot \left(\frac{d}{w_{\min}}\right)^{4\beta} \quad (11)$$

We now run the first step of the algorithm `Cluster` on X . We claim that the clustering obtained after the first step has very few classification errors. Let $\phi_r, r = 1, \dots, k$ be the k centers output by the first step of the algorithm `Cluster`. Lemma 5.1 implies that for each center θ_r , there exists a center ϕ_r satisfying

$$|\theta_r - \phi_r| \leq 20\sqrt{k} \cdot \frac{\|X - Z\|}{\sqrt{w_{\min} \cdot n}}.$$

Order the centers ϕ_r such that ϕ_r is closest to θ_r – equation (11) implies that the closest estimated centers to different θ_r are distinct. It also follows that for $r \neq s$

$$|\phi_r - \phi_s| \geq \frac{1}{2} \cdot \frac{\|X - Z\|}{\sqrt{n}} \cdot \left(\frac{d}{w_{\min}}\right)^{4\beta} \quad (12)$$

Lemma 7.5 *The number of points in S_r which are not assigned to ϕ_r after the first step of the algorithm is at most $\left(\frac{w_{\min}}{d}\right)^{2\beta} \cdot n$.*

Proof. We use the notation in Step 1 of algorithm `Cluster`. Suppose the statement of the lemma is not true. Then, in the k -means solution, at least $\left(\frac{w_{\min}}{d}\right)^{2\beta} \cdot n$ points in $\hat{X}_i, i \in S_r$ are assigned to a center at least $\frac{1}{2} \cdot \frac{\|X - Z\|}{\sqrt{n}} \cdot \left(\frac{d}{w_{\min}}\right)^{4\beta}$ distance away (using equation 12). But then the square of k -means clustering cost is much larger than $k \cdot \|X - Z\|^2$. ■

Now, we use the clustering given by the centers ϕ_r to partition the original set of points A – thus we have a clustering of these points where the number of *mis-classified* points from any cluster T_r is at most $\left(\frac{w_{\min}}{d}\right)^{2\beta} \cdot n$. Let S_r denote this clustering, where S_r corresponds to T_r . Let ν_r denote the center of S_r . We now argue that $|\nu_r - \mu_r|$ is very small.

Lemma 7.6 *For every $s, |\nu_s - \mu_s| \leq \frac{\|A - C\|}{\sqrt{n}}$.*

Proof. We use arguments similar to proof of Theorem 5.5. Let n_{rs}, μ_{rs} denote the number and mean respectively of $T_r \cap S_s$. Similarly, define n_{ss} and μ_{ss} as the size and mean of the points in $T_s \cap S_s$. We know that

$$|S_s| \nu_s = n_{ss} \mu_{ss} + \sum_{r \neq s} n_{rs} \mu_{rs}.$$

Theorem 5.4 implies that

$$|\mu_{rs} - \mu_s| \leq \frac{100 \cdot \|A - C\|}{\sqrt{n_{rs}}},$$

and Corollary 5.3 implies that

$$|\mu_{ss} - \mu_s| \leq \frac{\sqrt{|S_s| - n_{ss}}}{n_{ss}} \cdot \|A - C\|.$$

Now, proceeding as in the proof of Theorem 5.5, we get

$$\begin{aligned} |\mu_s - \nu_s| &\leq \frac{n_{ss}}{|S_s|} |\mu_{ss} - \mu_s| + \sum_{r \neq s} \frac{n_{rs}}{|S_s|} |\mu_{rs} - \mu_s| \\ &\leq \frac{400 \cdot \|A - C\|}{|T_s|} \cdot \left(\sum_{r \neq s} \sqrt{n_{rs}} \right) \end{aligned}$$

Using Lemma 7.5 now implies the result. ■

Starting from the centers ν_r , we run the second step of algorithm `Cluster`. Then, we have the analogue of Theorem 2.2 in this setting.

Theorem 7.7 *Suppose a mixture of distribution satisfies the conditions (8–10) above and at least $(1 - \varepsilon)$ fraction of sampled points satisfy the proximity condition. Then we can correctly classify all but $O(k^2 \varepsilon)$ fraction of the points.*

We now remove the assumption that $|\mu_r - \mu_s| \leq D_{\min} \cdot \left(\frac{d}{w_{\min}}\right)^\beta$ – let γ denote the latter quantity. We construct a graph $G = (V, E)$ as follows : V is the set of points $A \cup B$, and we join two points by an edge if the distance between them is at most γ/k . First observe that if $i, j \in T_r$, then they will be joined by an edge provided the following condition holds (using condition (9)) :

$$D_{\min} \cdot \sqrt{d} n^\alpha \text{polylog}(n) \leq D_{\min} \cdot \left(\frac{d}{w_{\min}}\right)^\beta,$$

and the same for A replaced by B above. This would hold if $\alpha < 0.1$ (recall that n is roughly $\left(\frac{d}{w_{\min}}\right)^{8(\beta+1)}$). Now consider the connected components of this graph. In each connected component, any two vertices are joined by a path of length at most k (because any two vertices from the same cluster T_r have an edge between them). So the distance between any two vertices from the same component is at most γ . Therefore the distance from the mean of two clusters in the same component is at most γ . Now, we can apply the arguments of this section to each component independently. This would, however, require us to know the number of clusters in each component of this graph. If we treat k as a constant, this is only constant number of choices. A better way is to modify the definition of X as follows : consider a point A_i . Let μ denote the mean of the points in the same component as A_i in the graph G . Then $X_{ij} = (A_i - \mu) \cdot (B_j - \mu)$ if A_i, B_j belong to the same component of G , L otherwise, where L is a large quantity. Now note that $\theta_r - \theta_s$ will still satisfy the statement of Lemma 7.1, because if they are from the same component in G , then it follows from the lemma, otherwise the distance between them is at least L . But Lemma 7.2 continues to hold without any change, and so rest of the arguments follow as they are.

7.1 Applications

We now give some applications of Theorem 7.7.

7.1.1 Learning Gaussian Distributions

Suppose we are given a mixture of Gaussian distribution F_1, \dots, F_k . Suppose the means satisfy the following separation condition for all $r, s, r \neq s$:

$$|\mu_r - \mu_s| \geq \Omega \left(\sigma k \cdot \log \frac{d}{w_{\min}} \right),$$

where σ denotes the maximum variance in any direction of the Gaussian distributions. Sample a set of $n = \text{poly} \left(\frac{d}{w_{\min}} \right)$ points. It is easy to check using Fact 6.1 that $\|A - C\|$ is $O(\sigma \sqrt{d} \cdot \log n)$. It is also easy to check that condition (10) is satisfied with $\alpha = 0$. Therefore, Theorem 7.7 implies the following.

Lemma 7.8 *Given a mixture of k Gaussians satisfying the separation condition above, we can correctly classify a set n samples, where $n = \text{poly} \left(\frac{d}{w_{\min}} \right)$.*

7.1.2 Learning Power Law Distributions

Consider a mixture of distributions F_1, \dots, F_k where each of the distributions F_r satisfies the following condition for every unit vector v :

$$P_{X \in F_r} [| (X - \mu_r) \cdot v | > \sigma t] \leq \frac{1}{t^\gamma} \quad (13)$$

where $\gamma \geq 2$ is a large enough constant. Let A be a set of n samples from the mixture. Suppose the means satisfy the following separation condition for every $r, s, r \neq s$:

$$|\mu_r - \mu_s| \geq \Omega \left(\sigma k \cdot \left(\log \frac{d}{w_{\min}} + \frac{1}{\varepsilon^\gamma} \right) \right).$$

First observe that since this is a special class of distributions considered in Section 6.3. So, one can again prove that $\frac{\|A - C\|}{\sqrt{n}}$ is $O(\sigma \cdot \sqrt{d} \cdot \text{polylog}(n))$. This is off from condition (8) by a factor of \sqrt{d} . But for large enough n , inequality (11) will continue to hold. Now let us try to bound $\max_i |A_i - C_i|$.

Claim 7.9 *With high probability,*

$$\max_i |A_i - C_i| \leq D_{\min} \cdot \sqrt{d} \cdot n^{\frac{2}{\gamma}} \cdot \text{polylog}(n).$$

Proof. Let e_1, \dots, e_d be orthonormal basis for the space. Then $| (A_i - C_i) \cdot e_l | \leq \sigma (nd)^{\frac{1}{\gamma}} \cdot \log(n)$ for all i with high probability. So, with high probability, for all i ,

$$|A_i - C_i| \leq D_{\min} \sqrt{dn}^{\frac{2}{\gamma}}.$$

■

Finally, we verify condition (10). Let v be a vector joining μ_r and μ_s . Then, $E \left[((A_i - C_i) \cdot v)^2 \right]$ is $O(\sigma^2)$ provided $\gamma \geq 2$. Now summing over all $A_i \in T_r$ and taking union bound for all k^2 choices for v proves that condition (10) is also satisfied. It is also easy to check that at least $1 - \varepsilon$ fraction of the points satisfy the proximity condition. So we have

Theorem 7.10 *Given a mixture of distributions where each distribution satisfies (13), we can cluster at least $1 - \varepsilon$ fraction of the points.*

References

- [ADK09] Ankit Aggarwal, Amit Deshpande, and Ravi Kannan. Adaptive sampling for k-means clustering. In *APPROX-RANDOM*, pages 15–28, 2009.
- [AK01] Sanjeev Arora and Ravi Kannan. Learning mixtures of arbitrary gaussians. In *ACM Symposium on Theory of Computing*, pages 247–257, 2001.
- [AM05] Dimitris Achlioptas and Frank McSherry. On spectral learning of mixtures of distributions. In *COLT*, pages 458–469, 2005.
- [AV06] David Arthur and Sergei Vassilvitskii. How slow is the k -means method? In *Symposium on Computational Geometry*, pages 144–153, 2006.
- [AV07] David Arthur and Sergei Vassilvitskii. k -means++: the advantages of careful seeding. In *SODA*, pages 1027–1035, 2007.
- [BV08] S. Charles Brubaker and Santosh Vempala. Isotropic pca and affine-invariant clustering. In *FOCS*, pages 551–560, 2008.
- [CR08] Kamalika Chaudhuri and Satish Rao. Beyond gaussians: Spectral methods for learning mixtures of heavy-tailed distributions. In *COLT*, pages 21–32, 2008.
- [Das99] Sanjoy Dasgupta. Learning mixtures of gaussians. In *IEEE Symposium on Foundations of Computer Science*, pages 634–644, 1999.
- [Das03] Sanjoy Dasgupta. How fast is k -means? In *COLT*, page 735, 2003.
- [DHKM07] Anirban Dasgupta, John E. Hopcroft, Ravi Kannan, and Pradipta Prometheus Mitra. Spectral clustering with limited independence. In *SODA*, pages 1036–1045, 2007.
- [DHKS05] Anirban Dasgupta, John E. Hopcroft, Jon M. Kleinberg, and Mark Sandler. On learning mixtures of heavy-tailed distributions. In *FOCS*, pages 491–500, 2005.
- [DLR77] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38, 1977.
- [DS07] Sanjoy Dasgupta and Leonard J. Schulman. A probabilistic analysis of em for mixtures of separated, spherical gaussians. *Journal of Machine Learning Research*, 8:203–226, 2007.
- [HPS05] Sariel Har-Peled and Bardia Sadri. How fast is the k -means method? In *SODA*, pages 877–885, 2005.
- [KSS10] Amit Kumar, Yogish Sabharwal, and Sandeep Sen. Linear-time approximation schemes for clustering problems in any dimensions. *J. ACM*, 57(2), 2010.
- [KSV08] Ravindran Kannan, Hadi Salmasian, and Santosh Vempala. The spectral method for general mixture models. *SIAM J. Comput.*, 38(3):1141–1156, 2008.
- [KV09] Ravi Kannan and Santosh Vempala. Spectral algorithms. *Foundations and Trends in Theoretical Computer Science*, 4(3-4):157–288, 2009.
- [Llo82] S. Lloyd. Least squares quantization in pcm. *Information Theory, IEEE Transactions on*, 28(2):129–137, 1982.

- [McS01] Frank McSherry. Spectral partitioning of random graphs. In *FOCS*, pages 529–537, 2001.
- [ORSS06] Rafail Ostrovsky, Yuval Rabani, Leonard J. Schulman, and Chaitanya Swamy. The effectiveness of lloyd-type methods for the k-means problem. In *FOCS*, pages 165–176, 2006.
- [VW04] Santosh Vempala and Grant Wang. A spectral algorithm for learning mixture models. *J. Comput. Syst. Sci.*, 68(4):841–860, 2004.